# CSE 486/586 Distributed Systems
## Domain Name System

Steve Ko
Computer Sciences and Engineering
University at Buffalo

CSE 486/586, Spring 2013

---

## Last Time

- Two multicast algorithms for total ordering
  – Sequencer
  – ISIS
- Multicast for causal ordering
  – Uses vector timestamps

CSE 486/586, Spring 2013                 2

---

## Review: Causally Ordered Multicast

- Each process keeps a vector clock.
  – Each counter represents the number of messages received from each of the other processes.
- When multicasting a message, the sender process increments its own counter and attaches its vector clock.
- Upon receiving a multicast message, the receiver process waits until it can preserve causal ordering:
  – It has delivered all the messages from the sender.
  – It has delivered all the messages that the sender had delivered before the multicast message.

CSE 486/586, Spring 2013                 3

---

## Review: Causal Ordering

Algorithm for group member $p_i$ ($i = 1, 2 \ldots, N$)

*On initialization*
$V_i^g[j] := 0 \ (j = 1, 2 \ldots, N);$

The number of group-g messages from process j that have been seen at process i so far

*To CO-multicast message m to group g*
$V_i^g[i] := V_i^g[i] + 1;$
$B\text{-multicast}(g, \langle V_i^g, m \rangle);$

*On B-deliver($\langle V_j^g, m \rangle$) from $p_j$, with g = group(m)*
place $\langle V_j^g, m \rangle$ in hold-back queue;
wait until $V_j^g[j] = V_i^g[j] + 1$ and $V_j^g[k] \le V_i^g[k]$ $(k \ne j);$
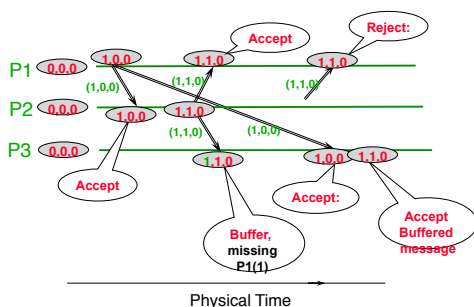CO-deliver m;    // after removing it from the hold-back queue
$V_i^g[j] := V_i^g[j] + 1;$

CSE 486/586, Spring 2013                 4

---

## Example: Causal Ordering Multicast



Physical Time

CSE 486/586, Spring 2013                 5

---

## Today's Question

- How do we organize the nodes in a distributed system?
- Up to the 90's
  – Prevalent architecture: client-server (or master-slave)
  – Unequal responsibilities
- Now
  – Emerged architecture: peer-to-peer
  – Equal responsibilities
- Studying an example client-server: DNS (today)
- Studying peer-to-peer as a paradigm (not just as a file-sharing application)
  – Learn the techniques and principles

CSE 486/586, Spring 2013                 6

C                                                                         1

## Separating Names and IP Addresses

- Names are easier (for us!) to remember
  - www.cnn.com vs. 64.236.16.20
- IP addresses can change underneath
  - Move www.cnn.com to 173.15.201.39
  - E.g., renumbering when changing providers
- Name could map to multiple IP addresses
  - www.cnn.com to multiple replicas of the Web site
- Map to different addresses in different places
  - Address of a nearby copy of the Web site
  - E.g., to reduce latency, or return different content
- Multiple names for the same address
  - E.g., aliases like ee.mit.edu and cs.mit.edu

## Two Kinds of Identifiers

- Host name (e.g., www.cnn.com)
  - Mnemonic name appreciated *by humans*
  - Provides little (if any) information about location
  - Hierarchical, variable # of alpha-numeric characters
- IP address (e.g., 64.236.16.20)
  - Numerical address appreciated *by routers*
  - Related to host's current location in the topology
  - Hierarchical name space of 32 bits

## Hierarchical Assignment Processes

- Host name: www.cse.buffalo.edu
  - Domain: registrar for each top-level domain (e.g., .edu)
  - Host name: local administrator assigns to each host
- IP addresses: 128.205.32.58
  - Prefixes: ICANN, regional Internet registries, and ISPs
  - Hosts: static configuration, or dynamic using DHCP

## Domain Name System (DNS)

Proposed in 1983 by Paul Mockapetris

## Overview: Domain Name System

- A client-server architecture
  - The server-side is still distributed for scalability.
  - But the servers are still a hierarchy of clients and servers
- Computer science concepts underlying DNS
  - Indirection: names in place of addresses
  - Hierarchy: in names, addresses, and servers
  - Caching: of mappings from names to/from addresses
- DNS software components
  - DNS resolvers
  - DNS servers
- DNS queries
  - Iterative queries
  - Recursive queries
- DNS caching based on time-to-live (TTL)

## Strawman Solution #1: Local File

- Original name to address mapping
  - Flat namespace
  - /etc/hosts
  - SRI kept main copy
  - Downloaded regularly
- Count of hosts was increasing: moving from a machine per domain to machine per user
  - Many more downloads
  - Many more updates

C

## Strawman Solution #2: Central Server

- Central server
  - One place where all mappings are stored
  - All queries go to the central server
- Many practical problems
  - Single point of failure
  - High traffic volume
  - Distant centralized database
  - Single point of update
  - Does not scale

Need a distributed, hierarchical collection of servers

## Domain Name System (DNS)

- Properties of DNS
  - Hierarchical name space divided into zones
  - Distributed over a collection of DNS servers
- Hierarchy of DNS servers
  - Root servers
  - Top-level domain (TLD) servers
  - Authoritative DNS servers
- Performing the translations
  - Local DNS servers
  - Resolver software

## DNS Root Servers

- 13 root servers (see http://www.root-servers.org/)
- Labeled A through M

A Verisign, Dulles, VA
C Cogent, Herndon, VA (also Los Angeles)
D U Maryland College Park, MD
G US DoD Vienna, VA
H ARL Aberdeen, MD
J Verisign, (11 locations)
E NASA Mt View, CA
F Internet Software C. Palo Alto, CA (and 17 other locations)
B USC-ISI Marina del Rey, CA
L ICANN Los Angeles, CA
K RIPE London (+ Amsterdam, Frankfurt)
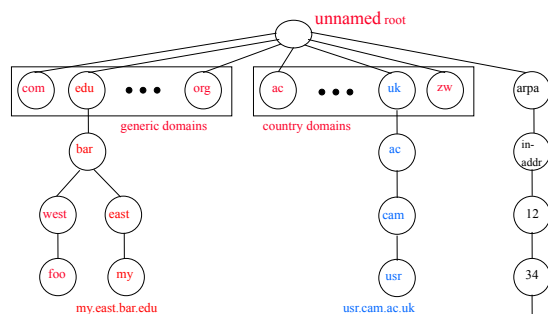I Autonomica, Stockholm (plus 3 other locations)
m WIDE Tokyo

## TLD and Authoritative DNS Servers

- Top-level domain (TLD) servers
  - Generic domains (e.g., com, org, edu)
  - Country domains (e.g., uk, fr, ca, jp)
  - Typically managed professionally
    » Network Solutions maintains servers for "com"
    » Educause maintains servers for "edu"
- Authoritative DNS servers
  - Provide public records for hosts at an organization
  - For the organization's servers (e.g., Web and mail)
  - Can be maintained locally or by a service provider

## Distributed Hierarchical Database

unnamed root

com  edu  • • •  org    ac  • • •  uk  zw    arpa

generic domains        country domains

bar                    ac              in-addr

west  east             cam             12

foo   my               usr             34

my.east.bar.edu        usr.cam.ac.uk   56

12.34.56.0/24

## Using DNS

- Local DNS server ("default name server")
  - Usually near the end hosts who use it
  - Local hosts configured with local server (e.g., /etc/resolv.conf) or learn the server via DHCP
- Client application
  - Extract server name (e.g., from the URL)
  - Do *gethostbyname()* to trigger resolver code
- Server application
  - Extract client IP address from socket
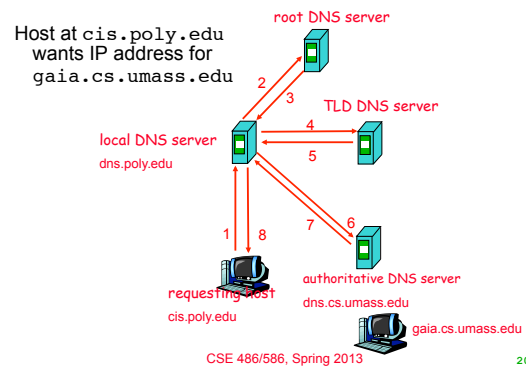  - Optional *gethostbyaddr()* to translate into name

## CSE 486/586 Administrivia

- Please start PA2 if you haven't.
- AWS codes are distributed on UBLearns.
  - Will post setup instructions.
- Practice problem set 1 & midterm example posted on the course website.
- Moving the midterm from Friday (3/8) to Wednesday (3/6)?
- Come talk to me!

---

## Example

Host at `cis.poly.edu` wants IP address for `gaia.cs.umass.edu`



root DNS server

TLD DNS server

local DNS server
dns.poly.edu

authoritative DNS server
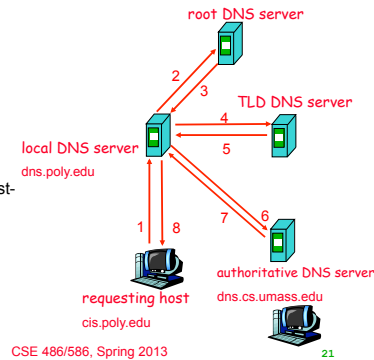dns.cs.umass.edu

requesting host
cis.poly.edu

gaia.cs.umass.edu

---

## Recursive vs. Iterative Queries

- Recursive query
  - Ask server to get answer for you
  - E.g., request 1 and response 8
- Iterative query
  - Ask server who to ask next
  - E.g., all other request-response pairs



root DNS server

TLD DNS server

local DNS server
dns.poly.edu

authoritative DNS server
dns.cs.umass.edu

requesting host
cis.poly.edu

---

## DNS Caching

- Performing all these queries take time
  - And all this before the actual communication takes place
  - E.g., 1-second latency before starting Web download
- Caching can substantially reduce overhead
  - The top-level servers very rarely change
  - Popular sites (e.g., www.cnn.com) visited often
  - Local DNS server often has the information cached
- How DNS caching works
  - DNS servers cache responses to queries
  - Responses include a "time to live" (TTL) field
  - Server deletes the cached entry after TTL expires

---

## Negative Caching

- Remember things that don't work
  - Misspellings like www.cnn.comm and www.cnnn.com
  - These can take a long time to fail the first time
  - Good to remember that they don't work
  - … so the failure takes less time the next time around

---

## DNS Resource Records

DNS: distributed db storing resource records (RR)

RR format: (name, value, type, ttl)

- Type=A
  - name is hostname
  - value is IP address

- Type=NS
  - `name` is domain (e.g. foo.com)
  - `value` is hostname of authoritative name server for this domain

- Type=CNAME
  - name is alias for some "canonical" (the real) name: www.ibm.com is really srveast.backup2.ibm.com
  - value is canonical name

- Type=MX
  - value is name of mailserver associated with name

C

4

## Reliability

- DNS servers are replicated
  - Name service available if at least one replica is up
  - Queries can be load balanced between replicas
- UDP used for queries
  - Need reliability: must implement this on top of UDP
- Try alternate servers on timeout
  - Exponential backoff when retrying same server
- Same identifier for all queries
  - Don't care which server responds

---

## Inserting Resource Records into DNS

- Example: just created startup "FooBar"
- Register foobar.com at Network Solutions
  - Provide registrar with names and IP addresses of your authoritative name server (primary and secondary)
  - Registrar inserts two RRs into the com TLD server:
    » (foobar.com, dns1.foobar.com, NS)
    » (dns1.foobar.com, 212.212.212.1, A)
- Put in authoritative server dns1.foobar.com
  - Type A record for www.foobar.com
  - Type MX record for foobar.com

- Play with "dig" on UNIX

---

```
$ dig nytimes.com ANY
; QUESTION SECTION:
;nytimes.com.                    IN      ANY

;; ANSWER SECTION:
nytimes.com.        267    IN     MX      100
  NYTIMES.COM.S7A1.PSMTP.com.
nytimes.com.        267    IN     MX      200
  NYTIMES.COM.S7A2.PSMTP.com.
nytimes.com.        267    IN     A       199.239.137.200
nytimes.com.        267    IN     A       199.239.136.200
nytimes.com.        267    IN     TXT "v=spf1 mx ptr
  ip4:199.239.138.0/24 include:alerts.wallst.com include:authsmtp.com
  ~all"
nytimes.com.        267    IN     SOA     ns1t.nytimes.com.
  root.ns1t.nytimes.com. 2009070102 1800 3600 604800 3600
nytimes.com.        267    IN     NS      nydns2.about.com.
nytimes.com.        267    IN     NS      ns1t.nytimes.com.
nytimes.com.        267    IN     NS      nydns1.about.com.

;; AUTHORITY SECTION:
nytimes.com.        267    IN     NS      nydns1.about.com.
nytimes.com.        267    IN     NS      ns1t.nytimes.com.
nytimes.com.        267    IN     NS      nydns2.about.com.

;; ADDITIONAL SECTION:
```

---

```
$ dig nytimes.com +norec @a.root-servers.net

;; ->>HEADER<<- opcode: QUERY, status: NOERROR, id: 53675
;; flags: qr; QUERY: 1, ANSWER: 0, AUTHORITY: 13, ADDITIONAL: 14

;; QUESTION SECTION:
;nytimes.com.                    IN      A

;; AUTHORITY SECTION:
com.            172800   IN     NS      K.GTLD-SERVERS.NET.
com.            172800   IN     NS      E.GTLD-SERVERS.NET.
com.            172800   IN     NS      D.GTLD-SERVERS.NET.
com.            172800   IN     NS      I.GTLD-SERVERS.NET.
com.            172800   IN     NS      C.GTLD-SERVERS.NET.

;; ADDITIONAL SECTION:
A.GTLD-SERVERS.NET.   172800   IN     A       192.5.6.30
A.GTLD-SERVERS.NET.   172800   IN     AAAA    2001:503:a83e::2:30
B.GTLD-SERVERS.NET.   172800   IN     A       192.33.14.30
B.GTLD-SERVERS.NET.   172800   IN     AAAA    2001:503:231d::2:30
```

---

```
$ dig nytimes.com +norec @k.gtld-servers.net

;; ->>HEADER<<- opcode: QUERY, status: NOERROR, id: 38385
;; flags: qr; QUERY: 1, ANSWER: 0, AUTHORITY: 3, ADDITIONAL: 3

;; QUESTION SECTION:
;nytimes.com.                    IN      A

;; AUTHORITY SECTION:
nytimes.com.        172800   IN     NS      ns1t.nytimes.com.
nytimes.com.        172800   IN     NS      nydns1.about.com.
nytimes.com.        172800   IN     NS      nydns2.about.com.

;; ADDITIONAL SECTION:
ns1t.nytimes.com.   172800   IN     A       199.239.137.15
nydns1.about.com.   172800   IN     A       207.241.145.24
nydns2.about.com.   172800   IN     A       207.241.145.25

;; Query time: 103 msec
;; SERVER: 192.52.178.30#53(192.52.178.30)
```

---

```
$ dig nytimes.com ANY +norec @ns1t.nytimes.com

;; ->>HEADER<<- opcode: QUERY, status: NOERROR, id: 39107
;; flags: qr aa; QUERY: 1, ANSWER: 13, AUTHORITY: 0, ADDITIONAL: 1

;; QUESTION SECTION:
;nytimes.com.                    IN      ANY

;; ANSWER SECTION:
nytimes.com.   300 IN    SOA     ns1t.nytimes.com.
  root.ns1t.nytimes.com. 2009070102 1800 3600 604800 3600
nytimes.com.   300 IN    MX      200 NYTIMES.COM.S7A2.PSMTP.com.
nytimes.com.   300 IN    MX      100 NYTIMES.COM.S7A1.PSMTP.com.
nytimes.com.   300 IN    NS      ns1t.nytimes.com.
nytimes.com.   300 IN    NS      nydns1.about.com.
nytimes.com.   300 IN    NS      nydns2.about.com.
nytimes.com.   300 IN    A       199.239.137.245
nytimes.com.   300 IN    A       199.239.136.200
nytimes.com.   300 IN    A       199.239.136.245
nytimes.com.   300 IN    TXT     "v=spf1...199.239.138.0/24
```
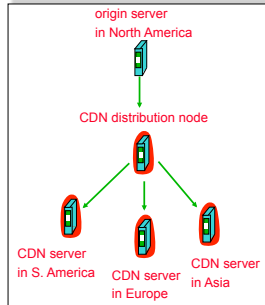
C

## Content Distribution Networks (CDNs)

- Content providers are CDN customers

<u>Content replication</u>
- CDN company installs thousands of servers throughout Internet
  - In large datacenters
  - Or, close to users
- CDN replicates customers' content
- When provider updates content, CDN updates servers



origin server
in North America

CDN distribution node

CDN server
in S. America

CDN server
in Europe

CDN server
in Asia

CSE 486/586, Spring 2013

---

## Content Distribution Networks

- Replicate content on many servers
- Challenges
  - How to replicate content
  - Where to replicate content
  - How to find replicated content
  - How to choose among replicas
  - How to direct clients towards a replica
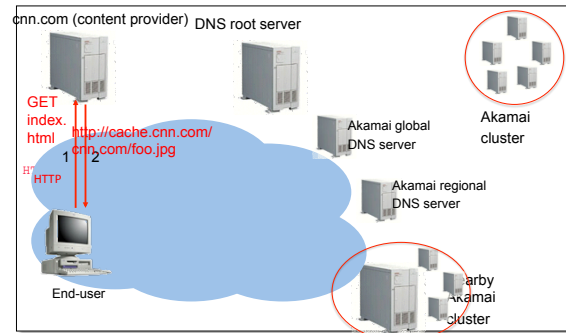
CSE 486/586, Spring 2013          32

---

## Server Selection

- Which server?
  - Lowest load: to balance load on servers
  - Best performance: to improve client performance
    » Based on what? Location? RTT? Throughput? Load?
  - Any alive node: to provide fault tolerance
- How to direct clients to a particular server?
  - As part of routing: anycast, cluster load balancer
  - As part of application: HTTP redirect
  - As part of naming: DNS
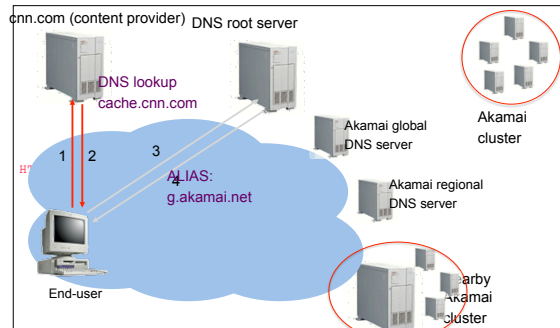
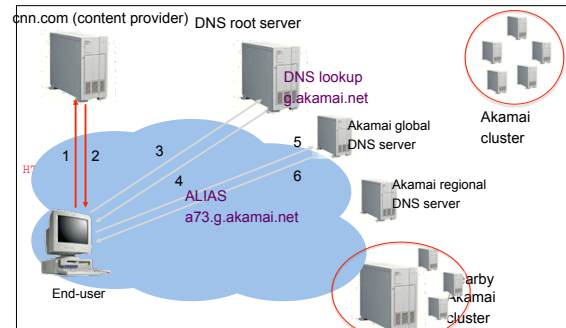CSE 486/586, Spring 2013          33

---

## How Akamai Works



cnn.com (content provider)   DNS root server

GET
index.
html

http://cache.cnn.com/
cnn.com/foo.jpg

1  2

HTTP

End-user

Akamai global
DNS server

Akamai
cluster

Akamai regional
DNS server

Nearby
Akamai
cluster

CSE 486/586, Spring 2013

---

## How Akamai Works



cnn.com (content provider)   DNS root server

DNS lookup
cache.cnn.com

1  2      3

4

ALIAS:
g.akamai.net

End-user

Akamai global
DNS server

Akamai
cluster

Akamai regional
DNS server

Nearby
Akamai
cluster

CSE 486/586, Spring 2013

---

## How Akamai Works



cnn.com (content provider)   DNS root server

DNS lookup
g.akamai.net

1  2      3            5

4            6

ALIAS
a73.g.akamai.net

End-user

Akamai global
DNS server

Akamai
cluster

Akamai regional
DNS server

Nearby
Akamai
cluster

CSE 486/586, Spring 2013

## How Akamai Works

cnn.com (content provider)  DNS root server

Akamai global
DNS server

Akamai
cluster

Akamai regional
DNS server

1  2     3
4          6
5          7
8
9

End-user

GET /foo.jpg
Host: cache.cnn.com

Nearby
Akamai
cluster

CSE 486/586, Spring 2013

## How Akamai Works

cnn.com (content provider)  DNS root server

GET foo.jpg

11
12

Akamai global
DNS server

Akamai
cluster

Akamai regional
DNS server

1  2     3
4          6
5          7
8
9

End-user

GET /foo.jpg
Host: cache.cnn.com

Nearby
Akamai
cluster

CSE 486/586, Spring 2013

## How Akamai Works

cnn.com (content provider)  DNS root server

11
12

Akamai global
DNS server

Akamai
cluster

Akamai regional
DNS server

1  2     3
4          6
5          7
8
9

End-user

10

Nearby
Akamai
cluster

CSE 486/586, Spring 2013

## Summary

- DNS as an example client-server architecture
- Why?
  - Names are easier (for us!) to remember
  - IP addresses can change underneath
  - Name could map to multiple IP addresses
  - Map to different addresses in different places
  - Multiple names for the same address
- Properties of DNS
  - Distributed over a collection of DNS servers
- Hierarchy of DNS servers
  - Root servers, top-level domain (TLD) servers, authoritative DNS servers

## Acknowledgements

- These slides contain material developed and copyrighted by Indranil Gupta (UIUC), Michael Freedman (Princeton), and Jennifer Rexford (Princeton).